

# Linear Algorithms for Online Multitask Classification

Giovanni Cavallanti\*

Nicolò Cesa-Bianchi†

Claudio Gentile‡

## Abstract

We design and analyze interacting online algorithms for multitask classification that perform better than independent learners whenever the tasks are related in a certain sense. We formalize task relatedness in different ways, and derive formal guarantees on the performance advantage provided by interaction. Our online analysis gives new stimulating insights into previously known co-regularization techniques, such as the multitask kernels and the margin correlation analysis for multiview learning. In the last part we apply our approach to spectral co-regularization: we introduce a natural matrix extension of the quasi-additive algorithm for classification and prove bounds depending on certain unitarily invariant norms of the matrix of task coefficients.

## 1 Introduction

A fundamental and fascinating problem in learning theory is the study of learning algorithms that influence each other. Although much is known about the behavior of individual strategies that learn a classification or regression task from examples, our understanding of interacting learning systems is still fairly limited. In this paper, we investigate this problem from the specific viewpoint of *multitask learning*, where each one of  $K > 1$  learners has to solve a different task (typically,  $K$  classification or  $K$  regression tasks). In particular, we focus on multitask binary classification, where learners are *online* linear classifiers (such as the Perceptron algorithm). Our goal is to design online interacting algorithms that perform better than independent learners whenever the tasks are related in a certain sense. We formalize task relatedness in different ways, and derive formal guarantees on the performance advantage provided by interaction.

Our analysis builds on ideas that have been developed in the context of statistical learning. In the statistical analysis of multitask learning (e.g., [2, 3, 4, 11, 24, 26]) the starting point is a regularized empirical loss functional or Tikhonov

functional —see, e.g., [10]. In the presence of several tasks, this functional is extended to allow for co-regularization among tasks. Roughly speaking, the co-regularization term forces the set of predictive functions for the  $K$  tasks to lie “close” to each other.

This co-regularization term is typically a squared norm in some Hilbert space of functions. We follow the approach pioneered by [11], where the  $K$  estimated solutions are linear functions parametrized by  $\mathbf{u} = (\mathbf{u}_1, \dots, \mathbf{u}_K) \in \mathbb{R}^{Kd}$  and the co-regularization is  $\mathbf{u}^\top A \mathbf{u}$ , where  $A$  is a positive definite matrix enforcing certain relations among tasks. The key observation in [11] is the following. Assume the instances of the multitask problem are of the form  $(\mathbf{x}_t, i_t)$ , where  $\mathbf{x}_t \in \mathbb{R}^d$  is an attribute vector and  $i_t \in \{1, \dots, K\}$  indicates the task  $\mathbf{x}_t$  refers to. Then one can reduce the  $K$  learning problems in  $\mathbb{R}^d$  to a single problem in  $\mathbb{R}^{Kd}$  by choosing a suitable embedding of the pairs  $(\mathbf{x}_t, i_t)$  into a common RKHS space  $\mathbb{R}^{Kd}$  with inner product  $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top A \mathbf{v}$ . This reduction allows us to solve a multitask learning problem by running any kernel-based single-task learning algorithm with the “multitask kernel” defined above. We build on this result by considering a natural online protocol for multitask linear classification. Within this protocol we analyze the performance of the Perceptron algorithm and some of its variants when run with a multitask kernel. Because such kernels are linear, we are not restricted to using kernel-based algorithms for efficiency reasons.

In Section 3 we consider the kernel Perceptron algorithm, and derive mistake bounds for the multitask kernels proposed in [11]. This reveals new insights into the role played by the regularizer matrix  $A$ . First, we see that the update in the kernel space defined by  $A$  factorizes in the “shared update” of  $K$  interacting Perceptrons each running in  $\mathbb{R}^d$ , thus providing a basic example of interactive online learning. Second, we exploit the simplicity of the mistake bound analysis to precisely quantify the performance advantage brought by the multitask approach over  $K$  independent online algorithms. In particular, in Subsections 3.2 and 3.3 we give examples where the mistake bound is used to guide the design of  $A$ . The first part of the paper is concluded with Sections 4 and 5, where we show multitask versions and mistake bound analyses for the second-order Perceptron algorithm of [7] and for the  $p$ -norm algorithm of [13, 14].

In the remaining sections of the paper, we depart from the approach of [11] to investigate the power of online learning when other forms of co-regularization are used. In Sec-

\*DSI, Università di Milano, Italy.

†DSI, Università di Milano, Italy.

‡DICOM, Università dell’Insubria, Italy. This is the corresponding author. Email: [claudio.gentile@uninsubria.it](mailto:claudio.gentile@uninsubria.it)

tion 6 we consider the case when instances belong to a space that is different for each task, and the similarity among tasks is measured by comparing their margin sequences (see, e.g., [6, 27]). We introduce and analyze a new multitask variant of the second-order Perceptron algorithm. The mistake bound that we prove is a margin-based version of the bound shown in Subsection 3.2 for the multitask Perceptron. Finally, in Section 7 we consider spectral co-regularization [4] for online multiview learning. Here diversity is penalized using a norm function defined on the  $d \times K$  matrix  $U = [\mathbf{u}_1, \dots, \mathbf{u}_K]$  of view vectors. In the spirit of [19], we interpret this penalization function as a potential defined over arbitrary matrices. We then define a natural extension of the quasi-additive algorithm of [14, 20] to a certain class of matrix norms, and provide a mistake bound analysis depending on the singular values of  $U$ . The results we obtain are similar to those in [28, 29, 30], though we are able to overcome some of the difficulties encountered therein via a careful study of matrix differentials.

In the next section, we introduce the basic online multitask protocol and define the multitask Perceptron algorithm. In order to keep the presentation as simple as possible, and to elucidate the interactive character of the updates, we delay the introduction of kernels until the proof of the mistake bound.

In our initial online protocol, at each time step the multitask learner receives a pair  $(\mathbf{x}_t, i_t)$ , where  $i_t$  is the task index for time  $t$  and  $\mathbf{x}_t$  is the instance for task  $i_t$ . Note that we view multitask learning as a sequential problem where at each time step the learner works on a single adversarially chosen task, rather than working simultaneously on all tasks (a similar protocol was investigated in [1] in the context of prediction with expert advice). One of the advantages of this approach is that, in most cases, the cost of running our multitask algorithms has a mild dependence on the number  $K$  of tasks.

We also remark that linear algorithms for online multitask learning have been studied in [9]. However, these results are sharply different from ours, as they do not depend on task relatedness.

## 2 Learning protocol and notation

There are  $K$  binary classification tasks indexed by  $1, \dots, K$ . At each time step  $t = 1, 2, \dots$  the learner receives a task index  $i_t \in \{1, \dots, K\}$  and the corresponding instance vector<sup>1</sup>  $\mathbf{x}_t \in \mathbb{R}^d$  (which we henceforth assume to be normalized,  $\|\mathbf{x}_t\| = 1$ ). Based on this information, the learner outputs a binary prediction  $\hat{y}_t \in \{-1, 1\}$  and then observes the correct label  $y_t \in \{-1, 1\}$  for task  $i_t$ . As in the standard worst-case online learning model, no assumptions are made on the mechanism generating the sequence  $(\mathbf{x}_t, y_t)_{t \geq 1}$ . Moreover, similarly to [1], the sequence of tasks  $i_t$  is also generated in an adversarial manner.

We compare the learner's performance to that of a reference predictor that is allowed to use a different linear classifier for each of the  $K$  tasks. In particular, we compare the

learner's mistake count to

$$\inf_{\mathbf{u}_1, \dots, \mathbf{u}_K \in \mathbb{R}^d} \sum_t \ell_t(\mathbf{u}_{i_t}) \quad (1)$$

where  $\ell_t(\mathbf{u}_{i_t}) = [1 - y_t \mathbf{u}_{i_t}^\top \mathbf{x}_t]_+$  is the hinge loss of the reference linear classifier (or *task vector*)  $\mathbf{u}_{i_t}$  at time  $t$ . Our goal is to design algorithms that make fewer mistakes than  $K$  independent learners when the tasks are related, and do not perform much worse than that when the tasks are completely unrelated. In the first part of the paper we use Euclidean distance to measure task relatedness. We say that the  $K$  tasks are related if there exist reference task vectors  $\mathbf{u}_1, \dots, \mathbf{u}_K \in \mathbb{R}^d$  having small pairwise distances  $\|\mathbf{u}_i - \mathbf{u}_j\|$ , and achieving a small cumulative hinge loss in the sense of (1). More general notions of relatedness are investigated in later sections.

## 3 The multitask Perceptron algorithm

We first introduce a simple multitask version of the Perceptron algorithm. This algorithm keeps a weight vector for each task and updates all weight vectors at each mistake using the Perceptron rule with different learning rates. More precisely, let  $\mathbf{w}_{i_t, t-1}$  be the weight vector associated with task  $i$  at the beginning of time step  $t$ . If we are forced (by the adversary) to predict on task  $i_t$ , and our prediction happens to be wrong, we update  $\mathbf{w}_{i_t, t-1}$  through the standard additive rule  $\mathbf{w}_{i_t, t} = \mathbf{w}_{i_t, t-1} + \eta y_t \mathbf{x}_t$  (where  $\eta > 0$  is a constant learning rate) but, at the same time, we perform a "half-update" on the remaining  $K - 1$  Perceptrons, i.e., we set  $\mathbf{w}_{j, t} = \mathbf{w}_{j, t-1} + \frac{\eta}{2} y_t \mathbf{x}_t$  for each  $j \neq i_t$ . This rule is based on the simple observation that, in the presence of related tasks, any update step that is good for one Perceptron should also be good for the others. Clearly, this rule keeps the weight vectors  $\mathbf{w}_{j, t}$ ,  $j = 1, \dots, K$ , always close to each other.

The above algorithm is a special case of the *multitask Perceptron algorithm* described below. This more general algorithm updates each weight vector  $\mathbf{w}_{j, t}$  through a learning rate which is an arbitrary positive definite function of the pair  $(j, i_t)$ . These learning rates are defined by a  $K \times K$  *interaction matrix*  $A$ .

The pseudocode for the multitask Perceptron algorithm using a generic interaction matrix  $A$  is given in Figure 1. At the beginning of each time step, the counter  $s$  stores the mistakes made so far (plus one). The (column) vector  $\phi_t \in \mathbb{R}^{Kd}$  denotes the *multitask instance* defined by

$$\phi_t^\top = \left( \underbrace{0, \dots, 0}_{d(i_t - 1) \text{ times}} \quad \mathbf{x}_t^\top \quad \underbrace{0, \dots, 0}_{d(K - i_t) \text{ times}} \right) \quad (2)$$

where  $\mathbf{x}_t \in \mathbb{R}^d$  is the instance vector for the current task  $i_t$ . (Note that  $\|\phi_t\| = 1$  since the instances  $\mathbf{x}_t$  are normalized.) The weights of the  $K$  Perceptrons are maintained in a compound vector  $\mathbf{w}_s^\top = (\mathbf{w}_{1, s}^\top, \dots, \mathbf{w}_{K, s}^\top)$ , with  $\mathbf{w}_{j, s} \in \mathbb{R}^d$  for all  $j$ . The algorithm predicts  $y_t$  through the sign  $\hat{y}_t$  of the  $i_t$ -th Perceptron's margin  $\mathbf{w}_{s-1}^\top \phi_t = \mathbf{w}_{i_t, s-1}^\top \mathbf{x}_t$ . Then, if prediction and true label disagree, the update rule becomes  $\mathbf{w}_s = \mathbf{w}_{s-1} + y_t (A \otimes I_d)^{-1} \phi_t$ , where  $\otimes$  denotes the Kro-

<sup>1</sup>Throughout this paper all vectors are assumed to be column vectors.

**Parameters:** Positive definite  $K \times K$  interaction matrix  $A$ .

**Initialization:**  $\mathbf{w}_0 = \mathbf{0} \in \mathbb{R}^{Kd}$ ,  $s = 1$ .

At each time  $t = 1, 2, \dots$  do the following:

1. Observe task number  $i_t \in \{1, \dots, K\}$  and the corresponding instance vector  $\mathbf{x}_t \in \mathbb{R}^d$ ;
2. Build the associated multitask instance  $\phi_t \in \mathbb{R}^{Kd}$ ;
3. Predict  $\hat{y}_t = \text{SGN}(\mathbf{w}_{s-1}^\top \phi_t) \in \{-1, +1\}$ ;
4. Get label  $y_t \in \{-1, +1\}$ ;
5. If  $\hat{y}_t \neq y_t$  then update:

$$\begin{aligned} \mathbf{w}_s &= \mathbf{w}_{s-1} + y_t (A \otimes I_d)^{-1} \phi_t \\ s &\leftarrow s + 1. \end{aligned}$$

Figure 1: The multitask Perceptron algorithm.

necker product between matrices<sup>2</sup> and  $I_d$  is the  $d \times d$  identity matrix. Since  $(A \otimes I_d)^{-1} = A^{-1} \otimes I_d$ , the above update is equivalent to the  $K$  task updates

$$\mathbf{w}_{j,s} \leftarrow \mathbf{w}_{j,s-1} + y_t A_{j,i_t}^{-1} \mathbf{x}_t \quad j = 1, \dots, K.$$

The algorithm is mistake driven, hence  $\mathbf{w}_{t-1}$  is updated (and is  $s$  increased) only when  $\hat{y}_t \neq y_t$ .

### 3.1 Pairwise distance interaction matrix

We now analyze the choice of  $A$  that corresponds to the updates  $\mathbf{w}_{i_t,s} \leftarrow \mathbf{w}_{i_t,s-1} + \eta y_t \mathbf{x}_t$  and  $\mathbf{w}_{j,s} \leftarrow \mathbf{w}_{j,s-1} + \frac{\eta}{2} y_t \mathbf{x}_t$  for  $j \neq i_t$  with  $\eta = 2/(K+1)$ . As it can be easily verified, this choice is given by

$$A = \begin{bmatrix} K & -1 & \dots & -1 \\ -1 & K & \dots & -1 \\ \dots & \dots & \dots & \dots \\ -1 & \dots & \dots & K \end{bmatrix} \quad (3)$$

with

$$A^{-1} = \frac{1}{K+1} \begin{bmatrix} 2 & 1 & \dots & 1 \\ 1 & 2 & \dots & 1 \\ \dots & \dots & \dots & \dots \\ 1 & \dots & \dots & 2 \end{bmatrix}.$$

In order to keep in with the notation just introduced, we equivalently specify an online multitask problem by the sequence  $(\phi_1, y_1), (\phi_2, y_2), \dots \in \mathbb{R}^{dK} \times \{-1, 1\}$  of multitask examples, where  $\phi_t$  is the multitask instance defined in (2). Moreover, given a sequence of multitask examples and reference task vectors  $\mathbf{u}_1, \dots, \mathbf{u}_K \in \mathbb{R}^d$ , we introduce the ‘‘compound’’ reference task vector  $\mathbf{u}^\top = (\mathbf{u}_1^\top, \dots, \mathbf{u}_K^\top) \in \mathbb{R}^{Kd}$  and write

$$\ell_t(\mathbf{u}) \stackrel{\text{def}}{=} [1 - y_t \mathbf{u}^\top \phi_t]_+ = [1 - y_t \mathbf{u}_{i_t}^\top \mathbf{x}_t]_+ = \ell_t(\mathbf{u}_{i_t}).$$

Finally, we use  $A_\otimes$  as a shorthand for  $A \otimes I_d$ , where  $d$  is understood from the context. We have the following result.

<sup>2</sup>The Kronecker or direct product between two matrices  $A = [a_{i,j}]$  and  $B$  of dimension  $m \times n$  and  $q \times r$ , respectively, is the block matrix of dimension  $m \times n$  whose block on row  $i$  and column  $j$  is the  $q \times r$  matrix  $a_{i,j} B$ .

**Theorem 1** The number of mistakes  $m$  made by the multitask Perceptron algorithm in Figure 1, run with interaction matrix (3) on any finite multitask sequence  $(\phi_1, y_1), (\phi_2, y_2), \dots \in \mathbb{R}^{Kd} \times \{-1, 1\}$ , satisfies

$$m \leq \inf_{\mathbf{u} \in \mathbb{R}^{Kd}} \left( \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u}) + \frac{2(\mathbf{u}^\top A_\otimes \mathbf{u})}{K+1} + \sqrt{\frac{2(\mathbf{u}^\top A_\otimes \mathbf{u})}{K+1} \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u})} \right),$$

where  $\mathcal{M}$  is the set of mistaken trial indices, and

$$\mathbf{u}^\top A_\otimes \mathbf{u} = \sum_{i=1}^K \|\mathbf{u}_i\|^2 + \sum_{1 \leq i < j \leq K} \|\mathbf{u}_i - \mathbf{u}_j\|^2.$$

**Remark** Note that when all tasks are equal, that is when  $\mathbf{u}_1 = \dots = \mathbf{u}_K$ , the bound of Theorem 1 becomes the standard Perceptron mistake bound (see, e.g., [7]). In the general case of distinct  $\mathbf{u}_i$  we have

$$\frac{\mathbf{u}^\top A_\otimes \mathbf{u}}{K+1} < \sum_{i=1}^K \|\mathbf{u}_i\|^2 - \frac{1}{K+1} \sum_{1 \leq i, j \leq K} \mathbf{u}_i^\top \mathbf{u}_j.$$

The sum of squares  $\sum_i \|\mathbf{u}_i\|^2$  is the mistake bound one can prove when learning  $K$  independent Perceptrons (under linear separability assumptions). On the other hand, highly correlated reference task vectors (i.e., large inner products  $\mathbf{u}_i^\top \mathbf{u}_j$ ) imply a large negative second term in the right-hand side of the above expression.

Theorem 1 is immediately proven by using the fact that the multitask Perceptron is a specific instance of the kernel Perceptron algorithm [12] using the linear kernel introduced in [11] (see also [15]). As mentioned in the introduction, this kernel is defined as follows: for any positive definite  $K \times K$  interaction matrix  $A$  introduce the  $Kd$ -dimensional reproducing kernel Hilbert space  $(\mathbb{R}^{Kd}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$  with inner product  $\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{H}} = \mathbf{u}^\top (A \otimes I_d) \mathbf{v}$ . Then define the kernel feature map  $\psi : \mathbb{R}^d \times \{1, \dots, K\} \rightarrow \mathcal{H}$  such that

$$\psi(\mathbf{x}_t, i_t) = (A \otimes I_d)^{-1} \phi_t. \quad (4)$$

The kernel used by the multitask Perceptron is thus defined by

$$\begin{aligned} K((\mathbf{x}_s, i_s), (\mathbf{x}_t, i_t)) &= \langle \psi(\mathbf{x}_s, i_s), \psi(\mathbf{x}_t, i_t) \rangle_{\mathcal{H}} \\ &= \phi_s^\top (A \otimes I_d)^{-1} \phi_t. \end{aligned} \quad (5)$$

**Proof of Theorem 1:** We use the following version of the kernel Perceptron bound (see, e.g., [7]),

$$\begin{aligned} m &\leq \sum_t \ell_t(f) + \|h\|_{\mathcal{H}}^2 \left( \max_t \|\psi(\mathbf{x}_t, i_t)\|_{\mathcal{H}}^2 \right) \\ &\quad + \|h\|_{\mathcal{H}} \sqrt{\left( \max_t \|\psi(\mathbf{x}_t, i_t)\|_{\mathcal{H}}^2 \right) \sum_t \ell_t(h)} \end{aligned}$$

where  $h$  is any function in the RKHS  $\mathcal{H}$  induced by the kernel. Let  $A_\otimes = A \otimes I_d$ . For the kernel (5) we have  $\|\mathbf{u}\|_{\mathcal{H}}^2 = \mathbf{u}^\top A_\otimes \mathbf{u}$  and  $\|\psi(\mathbf{x}_t, i_t)\|_{\mathcal{H}}^2 = \phi_t^\top A_\otimes^{-1} A_\otimes A_\otimes^{-1} \phi_t = \phi_t^\top A_\otimes^{-1} \phi_t = A_{i_t, i_t}^{-1}$ . Observing that  $A_{i_t, i_t}^{-1} = 2/(K+1)$  for the matrix  $A^{-1}$  defined in (3) concludes the proof. ■

### 3.2 A more general interaction matrix

In this section we slightly generalize the analysis of the previous section and consider an update rule of the form

$$\mathbf{w}_{j,s} = \mathbf{w}_{j,s-1} + \begin{cases} \frac{b+K}{(1+b)K} y_t \mathbf{x}_t & \text{if } j = i_t, \\ \frac{b}{(1+b)K} y_t \mathbf{x}_t & \text{otherwise,} \end{cases}$$

where  $b$  is a nonnegative parameter. The corresponding interaction matrix is given by

$$A = \frac{1}{K} \begin{bmatrix} a & -b & \dots & -b \\ -b & a & \dots & -b \\ \dots & \dots & \dots & \dots \\ -b & \dots & \dots & a \end{bmatrix}. \quad (6)$$

with  $a = K + b(K - 1)$ . It is immediate to see that the previous case (3) is recovered by choosing  $b = K$ . The inverse of (6) is

$$A^{-1} = \frac{1}{(1+b)K} \begin{bmatrix} b+K & b & \dots & b \\ b & b+K & \dots & b \\ \dots & \dots & \dots & \dots \\ b & \dots & \dots & b+K \end{bmatrix}.$$

When (6) is used in the multitask Perceptron algorithm, the proof of Theorem 1 can be adapted to prove the following result.

**Corollary 2** *The number of mistakes  $m$  made by the multitask Perceptron algorithm in Figure 1, run with interaction matrix (6) on any finite multitask sequence  $(\phi_1, y_1), (\phi_2, y_2), \dots \in \mathbb{R}^{Kd} \times \{-1, 1\}$ , satisfies*

$$m \leq \left( \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u}) + \frac{(b+K)}{(1+b)K} (\mathbf{u}^\top A_\otimes \mathbf{u}) + \sqrt{\frac{(b+K)}{(1+b)K} (\mathbf{u}^\top A_\otimes \mathbf{u}) \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u})} \right)$$

for any  $\mathbf{u} \in \mathbb{R}^{Kd}$ , where

$$\mathbf{u}^\top A_\otimes \mathbf{u} = \sum_{i=1}^K \|\mathbf{u}_i\|^2 + bK \text{VAR}[\mathbf{u}], \quad (7)$$

being  $\text{VAR}[\mathbf{u}] = \frac{1}{K} \sum_{i=1}^K \|\mathbf{u}_i - \bar{\mathbf{u}}\|^2$  the “variance”, of the task vectors, and  $\bar{\mathbf{u}}$  the centroid  $(\mathbf{u}_1 + \dots + \mathbf{u}_K)/K$ .

It is interesting to investigate how the above bound depends on the trade-off parameter  $b$ . The optimal value of  $b$  (requiring prior knowledge about the distribution of  $\mathbf{u}_1, \dots, \mathbf{u}_K$ ) is

$$b = \max \left\{ 0, \sqrt{\frac{K-1}{K} \frac{\|\bar{\mathbf{u}}\|^2}{\text{VAR}[\mathbf{u}]} - 1} \right\}.$$

Thus  $b$  grows large as the reference task vectors  $\mathbf{u}_i$  get close to their centroid  $\bar{\mathbf{u}}$  (i.e., as all  $\mathbf{u}_i$  get close to each other). Substituting this choice of  $b$  gives

$$\frac{(b+K)}{(1+b)K} (\mathbf{u}^\top A_\otimes \mathbf{u})$$

$$= \begin{cases} \|\mathbf{u}_1\|^2 + \dots + \|\mathbf{u}_K\|^2 & \text{if } b = 0, \\ \left( \|\bar{\mathbf{u}}\| + \sqrt{K-1} \sqrt{\text{VAR}[\mathbf{u}]} \right)^2 & \text{otherwise.} \end{cases}$$

When the variance  $\text{VAR}[\mathbf{u}]$  is large (compared to the squared centroid norm  $\|\bar{\mathbf{u}}\|^2$ ), then the optimal tuning of  $b$  is zero and the interaction matrix becomes the identity matrix, which amounts to running  $K$  independent Perceptron algorithms. On the other hand, when the optimal tuning of  $b$  is nonzero we learn  $K$  reference vectors, achieving a mistake bound equal to that of learning a *single* vector whose length is  $\|\bar{\mathbf{u}}\|$  plus  $\sqrt{K-1}$  times the standard deviation  $\sqrt{\text{VAR}[\mathbf{u}]}$ .

At the other extreme, if the variance  $\text{VAR}[\mathbf{u}]$  is zero (namely, when all tasks coincide) then the optimal  $b$  grows unbounded, and the quadratic term  $\frac{(b+K)}{(1+b)K} (\mathbf{u}^\top A_\otimes \mathbf{u})$  tends to the average square norm  $\frac{1}{K} \sum_{i=1}^K \|\mathbf{u}_i\|^2$ . In this case the multitask algorithm becomes essentially equivalent to an algorithm that, before learning starts, chooses one task at random and keeps referring all instance vectors  $\mathbf{x}_t$  to that task (somehow implementing the fact that now the information conveyed by  $i_t$  can be disregarded).

### 3.3 Encoding prior knowledge

We could also pick the interaction matrix  $A$  so to encode prior knowledge about tasks. For instance, suppose we know that only certain pairs of tasks are potentially related. We represent this knowledge in a standard way through an undirected graph  $G = (V, E)$ , where two vertices  $i$  and  $j$  are connected by an edge if and only if we believe task  $i$  and task  $j$  are related. A natural choice for  $A$  is then  $A = I + L$ , where  $L = [L_{i,j}]_{i,j=1}^K$  is the Laplacian of  $G$ , defined as

$$L_{i,j} = \begin{cases} d_i & \text{if } i = j, \\ -1 & \text{if } (i, j) \in E, \\ 0 & \text{otherwise,} \end{cases}$$

where  $d_i$  is the degree (number of incoming edges) of node  $i$ . If we now follow the proof of Theorem 1, which holds for any positive definite matrix  $A$ , we obtain the following result.

**Corollary 3** *The number of mistakes  $m$  made by the multitask Perceptron algorithm in Figure 1, run with interaction matrix  $I + L$  on any finite multitask sequence  $(\phi_1, y_1), (\phi_2, y_2), \dots \in \mathbb{R}^{Kd} \times \{-1, 1\}$ , satisfies*

$$m \leq \inf_{\mathbf{u} \in \mathbb{R}^{Kd}} \left( \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u}) + c_G \mathbf{u}^\top (I + L)_\otimes \mathbf{u} + \sqrt{c_G \mathbf{u}^\top (I + L)_\otimes \mathbf{u} \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u})} \right)$$

where

$$\mathbf{u}^\top (I + L)_\otimes \mathbf{u} = \sum_{i=1}^K \|\mathbf{u}_i\|^2 + \sum_{(i,j) \in E} \|\mathbf{u}_i - \mathbf{u}_j\|^2 \quad (8)$$

and  $c_G = \max_{i=1, \dots, K} \sum_{j=1}^K \frac{v_{j,i}^2}{1 + \lambda_j}$ . Here  $0 = \lambda_1 < \lambda_2 \leq \dots \leq \lambda_K$  are the eigenvalues of the positive semidef-

inite matrix  $L$ , and  $v_{j,i}$  denotes the  $i$ -th component<sup>3</sup> of the eigenvector  $\mathbf{v}_j$  of  $L$  associated with eigenvalue  $\lambda_j$ .

**Proof:** Following the proof of Theorem 1, we just need to bound

$$\max_{i=1,\dots,K} A_{i,i}^{-1} = \max_{i=1,\dots,K} (I + L)_{i,i}^{-1}.$$

If  $\mathbf{v}_1, \dots, \mathbf{v}_K$  are the eigenvectors of  $L$ , then

$$(I + L)^{-1} = \sum_{j=1}^K \frac{\mathbf{v}_j \mathbf{v}_j^\top}{1 + \lambda_j}$$

which concludes the proof.  $\blacksquare$

Ideally, we would like to have  $c_G = \mathcal{O}(1/K)$ . Clearly enough, if  $G$  is the clique on  $K$  vertices we expect to exactly recover the bound of Theorem 1. In fact, we can easily verify that the eigenvector  $\mathbf{v}_1$  associated with the zero eigenvalue  $\lambda_1$  is  $(K^{-1/2}, \dots, K^{-1/2})$ . Moreover, it is well known that all the remaining eigenvalues are equal to  $K$  (see, e.g., [16]). Therefore  $c_G = \frac{1}{K} + (1 - \frac{1}{K}) \frac{1}{K+1} = \frac{2}{K+1}$ . In the case of more general graphs  $G$ , we can bound  $c_G$  in terms of the smallest nonzero eigenvalue  $\lambda_2$ ,

$$c_G \leq \frac{1}{K} + \left(1 - \frac{1}{K}\right) \frac{1}{1 + \lambda_2}.$$

The value of  $\lambda_2$ , known as the algebraic connectivity of  $G$ , is 0 only when the graph is disconnected.  $\lambda_2$  is known for certain families of graphs. For instance, if  $G$  is a complete bipartite graph (i.e., if tasks can be divided in two disjoint subsets  $T_1$  and  $T_2$  such that every task in  $T_1$  is related to every task in  $T_2$  and for both  $i = 1, 2$  no two tasks in  $T_i$  are related), then it is known that  $\lambda_2 = \min\{|T_1|, |T_2|\}$ . We refer the reader to, e.g., [16] for further examples.

The advantage of using a graph  $G$  with significantly fewer edges than the clique is that the sum of pairwise distances in (8) will contain less than  $K(K-1)$  terms. On the other hand, this reduction is balanced by a larger coefficient  $c_G$  in front of  $\mathbf{u}^\top (I + L)_{\otimes} \mathbf{u}$ . This coefficient, in general, is related to the total number of edges in the graph (observe that the trace of  $L$  is exactly twice this total number).

## 4 The second-order extension

In this section we consider the second-order kernel Perceptron algorithm of [7] with the multitask kernel (5). The algorithm, which is described in Figure 2, maintains in its internal state a matrix  $S$  (initialized to the empty matrix) and a multitask Perceptron weight vector  $\mathbf{v}$  (initialized to the zero vector). Just like in Figure 1, we use the subscript  $s$  to denote the current number of mistakes plus one. Note that we have exploited the linearity of the kernel (5) to simplify the description of the algorithm. In particular, letting  $A_{\otimes} = A \otimes I_d$ , we have repeatedly used the fact that

$$\begin{aligned} \langle \psi(\mathbf{x}_s, i_s), \psi(\mathbf{x}_t, i_t) \rangle_{\mathcal{H}} &= \phi_s^\top A_{\otimes}^{-1} \phi_t \\ &= \left( A_{\otimes}^{-1/2} \phi_s \right)^\top \left( A_{\otimes}^{-1/2} \phi_t \right) \\ &= \tilde{\phi}_s^\top \tilde{\phi}_t, \end{aligned}$$

<sup>3</sup>Note that the orthonormality of the eigenvectors imply  $v_{1,i}^2 + \dots + v_{K,i}^2 = 1$  for all  $i$ .

**Parameters:** Positive definite  $K \times K$  interaction matrix  $A$ .  
**Initialization:**  $S_0 = \emptyset$ ,  $\mathbf{v}_0 = \mathbf{0} \in \mathbb{R}^{Kd}$ ,  $s = 1$ .

At each time  $t = 1, 2, \dots$  do the following:

1. Observe task number  $i_t \in \{1, \dots, K\}$  and the corresponding instance vector  $\mathbf{x}_t \in \mathbb{R}^d$ ;
2. Build the associated multitask instance  $\phi_t \in \mathbb{R}^{Kd}$  and compute  $\tilde{\phi}_t = (A \otimes I_d)^{-1/2} \phi_t$ ;
3. Predict  $\hat{y}_t = \text{SGN}(\mathbf{w}_{s-1}^\top \tilde{\phi}_t) \in \{-1, +1\}$ ,  
where  $\mathbf{w}_{s-1} = \left( I + S_{s-1} S_{s-1}^\top + \tilde{\phi}_t \tilde{\phi}_t^\top \right)^{-1} \mathbf{v}_{s-1}$ ;
4. Get the label  $y_t \in \{-1, +1\}$ ;
5. If  $\hat{y}_t \neq y_t$  then update:

$$\mathbf{v}_s = \mathbf{v}_{s-1} + y_t \tilde{\phi}_t, \quad S_s = \left[ S_{s-1} \mid \tilde{\phi}_t \right], \quad s \leftarrow s+1.$$

Figure 2: The second-order multitask Perceptron algorithm.

where  $\psi$  is the kernel feature map (4). The algorithm computes a tentative (inverse) matrix

$$\left( I + S_{s-1} S_{s-1}^\top + \tilde{\phi}_t \tilde{\phi}_t^\top \right)^{-1}.$$

Such a matrix is combined with the current Perceptron vector  $\mathbf{v}_{s-1}$  to predict the label  $y_t$ . If prediction  $\hat{y}_t$  and label  $y_t$  disagree both  $\mathbf{v}$  and  $S$  get updated (no update takes place otherwise). In particular, the new matrix  $S_s$  is augmented by padding with the current vector  $\tilde{\phi}_t$ . Since supports are shared, the computational cost of an update is not significantly larger than that for learning a single-task (see Section 4.1).

**Theorem 4** *The number of mistakes  $m$  made by the second-order multitask Perceptron of Figure 2, run with any positive definite interaction matrix  $A$ , on any finite multitask sequence  $(\phi_1, y_1), (\phi_2, y_2), \dots \in \mathbb{R}^{Kd} \times \{-1, +1\}$ , satisfies, for all  $\mathbf{u} \in \mathbb{R}^{Kd}$ ,*

$$m \leq \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u}) + \sqrt{\left( \mathbf{u}^\top (A \otimes I_d) \mathbf{u} + \sum_{t \in \mathcal{M}} (\mathbf{u}_{i_t}^\top \mathbf{x}_t)^2 \right) \sum_{j=1}^m \ln(1 + \lambda_j)}$$

where  $\mathcal{M}$  is the sequence of mistaken trial indices and  $\lambda_1, \dots, \lambda_m$  are the eigenvalues of the  $m \times m$  matrix of elements  $\mathbf{x}_s^\top A_{i_s, i_t}^{-1} \mathbf{x}_t$ , where  $s, t \in \mathcal{M}$ .

**Proof:** Recall the mistake bound for the second-order kernel Perceptron algorithm [7]:

$$m \leq \sqrt{\left( \|h\|_{\mathcal{H}}^2 + \sum_{t \in \mathcal{M}} h(\mathbf{x}_t)^2 \right) \sum_{j=1}^m \ln(1 + \lambda_j)}$$

where  $\lambda_1, \dots, \lambda_m$  are the eigenvalues of the  $m \times m$  kernel Gram (sub)matrix including only time steps in  $\mathcal{M}$ . When the

kernel is (5) with feature map  $f$  we have  $\|\mathbf{u}\|_{\mathcal{H}}^2 = \mathbf{u}^\top A_\otimes \mathbf{u}^\top$  and  $\langle \mathbf{u}^\top, \psi(\mathbf{x}_t, i_t) \rangle^2 = (\mathbf{u}^\top A_\otimes A_\otimes^{-1} \phi_t)^2 = (\mathbf{u}_t^\top \mathbf{x}_t)^2$ . Finally, the kernel Gram matrix is  $K(\psi(\mathbf{x}_s, i_s), \psi(\mathbf{x}_t, i_t)) = \phi_s^\top A_\otimes^{-1} \phi_t = \mathbf{x}_s^\top A_{i_s, i_t}^{-1} \mathbf{x}_t$ . This concludes the proof.  $\blacksquare$

Again, this bound should be compared to the one obtained when learning  $K$  independent tasks. As in the first-order algorithm, we have the complexity term  $\mathbf{u}^\top (A \otimes I_d) \mathbf{u}$ . In this case, however, the interaction matrix  $A$  also plays a role in the scale of the eigenvalues of the resulting multitask Gram matrix. Roughly speaking, we gain a factor  $K$  from  $\mathbf{u}^\top A_\otimes^{-1} \mathbf{u}$  (according to the arguments in Section 3). In addition, however, we gain a further factor  $K$ , since the trace of the multitask Gram matrix  $[\phi_s^\top A_\otimes^{-1} \phi_t]_{s,t \in \mathcal{M}} = [\mathbf{x}_s^\top A_{i_s, i_t}^{-1} \mathbf{x}_t]_{s,t \in \mathcal{M}}$  is about  $1/K$  times the trace of the original Gram matrix  $[\mathbf{x}_s^\top \mathbf{x}_t]_{s,t \in \mathcal{M}}$ . Since both factors are under the square root, the resulting gain over the  $K$  independent task bound is about  $K$ .

#### 4.1 Implementation in dual variables

It is easy to see that the second-order multitask Perceptron can be run in dual variables by maintaining  $K$  classifiers that share the same set of support vectors. This allows an efficient implementation that does not impose any significant overhead with respect to the corresponding single-task version. Specifically, given some interaction matrix  $A$ , the margin at time  $t$  is computed as (see [7, Theorem 3.3])

$$\begin{aligned} \mathbf{w}_{s-1}^\top \tilde{\phi}_t &= \mathbf{v}_{s-1}^\top \left( I + S_{s-1} S_{s-1}^\top + \tilde{\phi}_t \tilde{\phi}_t^\top \right)^{-1} \tilde{\phi}_t \\ &= \mathbf{y}_s^\top \left( I + S_s^\top S_s \right)^{-1} S_s^\top \tilde{\phi}_t, \end{aligned} \quad (9)$$

where  $\mathbf{y}_s$  is the  $s$ -dimensional vector whose first  $s-1$  components are the labels  $y_i$  where the algorithm has made a mistake up to time  $t-1$ , and the last component is 0.

First, note that replacing  $I + S_s^\top S_s$  with  $I + S_{s-1}^\top S_{s-1}$  in (9) does not change the sign of the prediction. The margin at time  $t$  can then be computed by calculating the scalar product between  $S_s^\top \tilde{\phi}_t$  and  $\mathbf{y}_s^\top (I + S_{s-1}^\top S_{s-1})^{-1}$ . Now, each entry of the vector  $S_s^\top \tilde{\phi}_t$  is of the form  $A_{j,i_t}^{-1} \mathbf{x}_j^\top \mathbf{x}_t$ , and thus computing  $S_s^\top \tilde{\phi}_t$  requires  $\mathcal{O}(s)$  inner products so that, overall, the prediction step requires  $\mathcal{O}(s)$  scalar multiplications and  $\mathcal{O}(s)$  inner products (independent of the number of tasks  $K$ ).

On the other hand, the update step involves the computation of the vector  $\mathbf{y}_s^\top (I + S_s^\top S_s)^{-1}$ . For the matrix update we can write

$$I_s + S_s^\top S_s = \begin{bmatrix} I_{s-1} + S_{s-1}^\top S_{s-1} & S_{s-1}^\top \tilde{\phi}_t \\ \tilde{\phi}_t^\top S_{s-1} & 1 + \tilde{\phi}_t^\top \tilde{\phi}_t \end{bmatrix}.$$

Using standard facts about the inverse of partitioned matrices (e.g., [17, Ch. 0]), one can see that the inverse of matrix  $I_s + S_s^\top S_s$  can be computed from the inverse of  $I_{s-1} + S_{s-1}^\top S_{s-1}$  with  $\mathcal{O}(s)$  extra inner products (again, independent of  $K$ ) and  $\mathcal{O}(s^2)$  additional scalar multiplications.

## 5 The $p$ -norm extension

We now extend our multitask results to the  $p$ -norm Perceptron algorithm of [14, 13]. As before, when the tasks are all equal we want to recover the bound of the single-task algorithm, and when the task vectors are different we want the mistake bound to increase according to a function that penalizes task diversity according to their  $p$ -norm distance.

We develop our  $p$ -norm multitask analysis for the specific choices of  $p = 2 \ln d$  (or  $p = 2 \ln K$  when  $d \leq K$ ) and for the pairwise distance matrix (3). It is well known that for  $p = 2 \ln d$  the mistake bound of the single-task  $p$ -norm Perceptron is essentially equivalent to the one of the zero-threshold Winnow algorithm of [22]. We now see that this property is preserved in the multitask extension.

We start with the following slightly more general algorithm based on arbitrary norms. Later, we specialize it to  $p$ -norms. The *quasi-additive multitask algorithm* of [20, 14] is defined for any norm  $\|\cdot\|$  over  $\mathbb{R}^{Kd}$ . Initially,  $\mathbf{w}_0 = \mathbf{0} \in \mathbb{R}^{Kd}$ . If  $s-1$  mistakes have been made in the first  $t-1$  time steps, then the prediction at time  $t$  is  $\text{SGN}(\mathbf{w}_{s-1}^\top \phi_t)$ . If a mistake occurs at time  $t$ , then  $\mathbf{w}_{s-1}$  is updated with the rule  $\mathbf{w}_s = \nabla_{\frac{1}{2}} \|\mathbf{v}_s\|^2$ , where the *primal weight*  $\mathbf{v}_s \in \mathbb{R}^{Kd}$  is updated using the multitask Perceptron rule,  $\mathbf{v}_0 = \mathbf{0} \in \mathbb{R}^{Kd}$  and  $\mathbf{v}_s = \mathbf{v}_{s-1} + y_t A_\otimes^{-1} \phi_t$  for an arbitrary positive definite interaction matrix  $A$ .

This can be analyzed using the following technique (see [8] for details). Let  $\mathbf{v}_m$  be the primal weight after any number  $m$  of mistakes. Then, by Taylor expanding  $\frac{1}{2} \|\mathbf{v}_s\|^2$  around  $\mathbf{v}_{s-1}$  for each  $s = 1, \dots, m$ , and using the fact  $y_t \mathbf{w}_{s-1}^\top \phi_t \leq 0$  whenever a mistake occurs at step  $t$ , we get

$$\frac{1}{2} \|\mathbf{v}_m\|^2 \leq \sum_{s=1}^m D(\mathbf{v}_s \| \mathbf{v}_{s-1}) \quad (10)$$

$D(\mathbf{v}_s \| \mathbf{v}_{s-1}) = \frac{1}{2} (\|\mathbf{v}_s\|^2 - \|\mathbf{v}_{s-1}\|^2) - y_t \mathbf{w}_{s-1}^\top \mathbf{x}_t$  is a so-called Bregman divergence; i.e., the error term in the first-order Taylor expansion of  $\frac{1}{2} \|\cdot\|^2$  around vector  $\mathbf{v}_{s-1}$ , at vector  $\mathbf{v}_s$ .

Fix any  $\mathbf{u} \in \mathbb{R}^{Kd}$ . Using the convex inequality for norms  $\mathbf{u}^\top \mathbf{v} \leq \|\mathbf{u}\| \|\mathbf{v}\|_*$  where  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$  (see, e.g., [25, page 131]), and using the fact  $\mathbf{u}^\top A_\otimes \mathbf{v}_s = \mathbf{u}^\top A_\otimes \mathbf{v}_{s-1} + y_t \mathbf{u}^\top \phi_t \geq \mathbf{u}^\top A_\otimes \mathbf{v}_{s-1} + 1 - \ell_t(\mathbf{u})$ , one then obtains

$$\|\mathbf{v}_m\| \geq \frac{\mathbf{u}^\top A_\otimes \mathbf{v}_m}{\|A_\otimes \mathbf{u}\|_*} \geq \frac{m - \sum_t \ell_t(\mathbf{u})}{\|A_\otimes \mathbf{u}\|_*}. \quad (11)$$

Combining (10) with (11) and solving for  $m$  gives

$$m \leq \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u}) + \|A_\otimes \mathbf{u}\|_* \sqrt{2 \sum_{s=1}^m D(\mathbf{v}_s \| \mathbf{v}_{s-1})}. \quad (12)$$

We obtain our multitask version of the  $p$ -norm Perceptron when  $\|\mathbf{u}\| = \|\mathbf{u}\|_p = (|u_1|^p + |u_2|^p + \dots)^{1/p}$ . In particular, we focus our analysis on the choice  $p = 2 \ln \max\{K, d\}$ , which gives mistake bounds in the dual norms  $\|\mathbf{u}\|_1$  and  $\|\mathbf{x}_t\|_\infty$ , and on the pairwise distance matrix (3).

Using the analysis in [8] we obtain, for  $t_s = t$ ,

$$D(\mathbf{v}_s \| \mathbf{v}_{s-1}) \leq \frac{p-1}{2} \|A_\otimes^{-1} \phi_t\|_p^2 = \frac{p-1}{2} \|\mathbf{x}_t\|_p^2 \|A_{\downarrow i_t}^{-1}\|_p^2$$

where  $A_{\downarrow i_t}^{-1}$  is the  $i_t$ -th column of  $A^{-1}$ . If we now use  $p = 2 \ln \max\{K, d\}$ , then  $\|\mathbf{x}_t\|_p^2 \leq e \|\mathbf{x}_t\|_\infty^2$  and

$$\left\|A_{\downarrow i_t}^{-1}\right\|_p^2 \leq e \left\|A_{\downarrow i_t}^{-1}\right\|_\infty^2 = e (A_{i_t, i_t}^{-1})^2 = \frac{4e}{(K+1)^2}.$$

We now turn to the computation of the dual norm  $\|A_\otimes \mathbf{u}\|_q$ , where  $q = p/(p-1)$  is the dual coefficient of  $p$ . We find that

$$\|A_\otimes \mathbf{u}\|_q^2 \leq \|A_\otimes \mathbf{u}\|_1^2 = \left( \sum_{i=1}^K \left\| \mathbf{u}_i + \sum_{j \neq i} (\mathbf{u}_i - \mathbf{u}_j) \right\|_1 \right)^2.$$

Plugging back into (12) gives the following theorem.

**Theorem 5** *The number of mistakes  $m$  made by the  $p$ -norm multitask Perceptron, run with the pairwise distance matrix (3) and  $p = 2 \ln \max\{K, d\}$ , on any finite multitask sequence  $(\phi_1, y_1), (\phi_2, y_2), \dots \in \mathbb{R}^{Kd} \times \{-1, 1\}$ , satisfies, for all  $\mathbf{u} \in \mathbb{R}^{Kd}$ ,*

$$m \leq \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u}) + H + \sqrt{2H \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u})}$$

where  $H$  is equal to

$$\frac{4e^2 \ln \max\{K, d\}}{(K+1)^2} X_\infty^2 \left( \sum_{i=1}^K \left\| \mathbf{u}_i + \sum_{j \neq i} (\mathbf{u}_i - \mathbf{u}_j) \right\|_1 \right)^2.$$

and  $X_\infty = \max_{t \in \mathcal{M}} \|\mathbf{x}_t\|_\infty$ .

**Remark** When all tasks are equal,  $\mathbf{u}_1 = \dots = \mathbf{u}_K$ , the coefficient  $H$  in the bound of Theorem 5 becomes

$$(4e^2 \ln \max\{K, d\}) \left( \max_t \|\mathbf{x}_t\|_\infty \right)^2 \|\mathbf{u}_i\|_1^2.$$

If  $K \leq d$  this bound is equivalent (apart from constant factors) to the mistake bound for the single-task zero-threshold Winnow algorithm of [22].

## 6 Learning tasks in heterogeneous spaces

In this section we slightly deviate from the approach followed so far. We consider the case when the  $K$  task vectors  $\mathbf{u}_i$  may live in different spaces:  $\mathbf{u}_i \in \mathbb{R}^{d_i}$ ,  $i = 1, \dots, K$ . This is a plausible assumption when attributes associated with different tasks have a completely different meaning. In such a case, the correlation among tasks is naturally measured through the task margins  $\mathbf{u}_i^\top \mathbf{x}$  (or *views*)—see, e.g., the previous work of [6, 27] for a similar approach in the context of semi-supervised learning. In order to allow for the views to interact in a meaningful way, we slightly modify the learning protocol of Section 2. We now assume that, at each time  $t$ , we receive the adversarial choice of task  $i_t$  together with *all* instance vectors  $\mathbf{x}_{i,t} \in \mathbb{R}^{d_i}$  for  $i = 1, \dots, K$ . The co-regularization terms motivating our algorithm are proportional to the distance between the margin  $\mathbf{u}_{i_t}^\top \mathbf{x}_{i_t,t}$  of task  $i_t$  and the average margin  $\frac{1}{K} \sum_{j=1}^K \mathbf{u}_j^\top \mathbf{x}_{j,t}$  of all tasks:

$$bK \left( \mathbf{u}_{i_t}^\top \mathbf{x}_{i_t,t} - \frac{1}{K} \sum_{j=1}^K \mathbf{u}_j^\top \mathbf{x}_{j,t} \right)^2, \quad (13)$$

where  $b$  is a positive constant. We add the above terms up to time  $t$ , resulting in a cumulative regularization term

$$bK \sum_{s \in \mathcal{M}_t} \left( \mathbf{u}_{i_s}^\top \mathbf{x}_{i_s,s} - \frac{1}{K} \sum_{j=1}^K \mathbf{u}_j^\top \mathbf{x}_{j,s} \right)^2,$$

where  $\mathcal{M}_t$  is the set of mistaken trials up to time  $t$ . Thus, in trial  $t$  our prior knowledge about task relatedness is encoded as a (positive) correlation among the  $K$  margin sequences  $(\mathbf{u}_1^\top \mathbf{x}_{j,1}, \mathbf{u}_j^\top \mathbf{x}_{j,2}, \dots, \mathbf{u}_j^\top \mathbf{x}_{j,t})$ ,  $j = 1, \dots, K$ .

The algorithm of Figure 3 is a natural multitask predictor operating with the above view-based regularization criterion. We call this algorithm the multiview-based multitask Perceptron algorithm, or MMPERC for brevity. MMPERC can be viewed as a variant of the second-order Perceptron using the cumulative covariance matrix of the past margin vectors in order to suitably transform instances.

MMPERC has a constant tradeoff parameter  $b > 0$  (playing the same role as the one in Section 3.2), and maintains in its internal state a multitask matrix  $A$  (initialized to the identity matrix  $I$ ) and a Perceptron multitask weight vector  $\mathbf{v}$  (initialized to the zero vector). The subscript  $s$  plays the same role as in the previous algorithms. Unlike the algorithms in previous sections, MMPERC observes the task number  $i_t$  and the multitask instance  $\Phi_t^\top = (\mathbf{x}_{1,t}^\top, \mathbf{x}_{2,t}^\top, \dots, \mathbf{x}_{K,t}^\top)$  made up of the instance vectors of all tasks. Then MMPERC computes a tentative matrix  $A'_t$  to be used for prediction. Matrix  $A'_t$  is obtained by adding the rank-one positive semidefinite matrix  $M_t$  to the previous matrix  $A_{s-1}$ . Here  $\phi_{j,t}$  is the  $(d_1 + \dots + d_K)$ -dimensional vector

$$\phi_{j,t}^\top = \left( \underbrace{0, \dots, 0}_{d_1 + \dots + d_{j-1} \text{ times}} \quad \mathbf{x}_{j,t}^\top \quad \underbrace{0, \dots, 0}_{d_{j+1} + \dots + d_K \text{ times}} \right)$$

for  $j = 1, \dots, K$ . Observe that  $M_t$  has been set so as to make the quadratic form  $\mathbf{u}^\top M_t \mathbf{u}$  coincide with the regularization term (13). Similarly to the algorithms of previous sections, the tentative matrix  $A'_t$  and the current Perceptron vector  $\mathbf{v}_{s-1}$  are used for predicting the true label  $y_t$ . If prediction  $\hat{y}_t$  and label  $y_t$  disagree both  $\mathbf{v}$  and  $A$  get updated. In particular,  $A_s$  is set to the tentative matrix  $A'_t$ .

In this protocol we call example the triple  $(i_t, \Phi_t, y_t)$ . Like the results contained in the previous sections, our analysis will provide a multitask bound on the number of prediction mistakes which is comparable to the one obtained by a single task plus a penalization term due to task relatedness. However, though this algorithm is a second-order prediction method, we only give a first-order analysis that disregards the eigenstructure of the data. This is due to the technical difficulty of handling a time-varying matrix  $A$  that in trial  $t$  includes all instance vectors  $\mathbf{x}_{1,t}, \dots, \mathbf{x}_{K,t}$ .

**Theorem 6** *The number of mistakes  $m$  made by the algorithm in Figure 3, run on any multitask sequence  $(i_1, \Phi_1, y_1), (i_2, \Phi_2, y_2), \dots$  satisfies, for all  $\mathbf{u}^\top =$*

**Parameters:**  $b > 0$ .

**Initialization:**  $A_0 = I$ ,  $\mathbf{v}_0 = \mathbf{0} \in \mathbb{R}^{d_1 + \dots + d_K}$ ,  $s = 1$ .

At each time  $t = 1, 2, \dots$  do the following:

1. Observe task number  $i_t \in \{1, \dots, K\}$ ;

2. Observe multitask instance vector

$$\Phi_t^\top = (\mathbf{x}_{1,t}^\top, \dots, \mathbf{x}_{K,t}^\top) \in \mathbb{R}^{d_1 + \dots + d_K};$$

3. Build the associated multitask instance  $\phi_{i_t,t}$ ;

4. Set  $A'_t = A_{s-1} + M_t$  where

$$M_t = bK \left( \phi_{i_t,t} - \frac{\Phi_t}{K} \right) \left( \phi_{i_t,t} - \frac{\Phi_t}{K} \right)^\top;$$

5. Predict  $\hat{y}_t = \text{SGN}(\mathbf{v}_{s-1}^\top (A'_t)^{-1} \phi_{i_t,t}) \in \{-1, +1\}$ ;

6. Get label  $y_t \in \{-1, +1\}$ ;

7. If  $\hat{y}_t \neq y_t$  then update:

$$\mathbf{v}_s = \mathbf{v}_{s-1} + y_t \phi_{i_t,t}, \quad A_s = A'_t, \quad s \leftarrow s + 1.$$

Figure 3: The multiview-based multitask Perceptron algorithm (MMPERC).

$$(\mathbf{u}_1^\top, \dots, \mathbf{u}_K^\top),$$

$$\begin{aligned} m \leq \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u}) + \frac{K(b+1) - b}{bK(K-1) + K} (\mathbf{u}^\top A_m \mathbf{u}) \\ + \sqrt{\frac{K(b+1) - b}{bK(K-1) + K} (\mathbf{u}^\top A_m \mathbf{u}) \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u})}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{u}^\top A_m \mathbf{u} \\ = \sum_{i=1}^K \|\mathbf{u}_i\|^2 + bK \sum_{t \in \mathcal{M}} \left( \mathbf{u}_{i_t}^\top \mathbf{x}_{i_t,t} - \frac{1}{K} \sum_{j=1}^K \mathbf{u}_j^\top \mathbf{x}_{j,t} \right)^2, \end{aligned}$$

being  $\mathcal{M}$  the set of mistaken trials.

**Remark** It is the factor

$$\frac{K(b+1) - b}{bK(K-1) + K} (\mathbf{u}^\top A_m \mathbf{u}) \quad (14)$$

that quantifies the relatedness among tasks (the leading constant  $bK$  in the second term of  $\mathbf{u}^\top A_m \mathbf{u}$  is needed for scaling purposes). Note that the notion of relatedness provided by  $\mathbf{u}^\top A_m \mathbf{u}$  is analogous to the one used in Section 3.2. As suggested in Section 3.3, other measures of similarity are possible.

The parameter  $b$  allows for a limited trade-off between  $\sum_{i=1}^K \|\mathbf{u}_i\|^2$  and the cumulative “margin deviation”

$$\sum_{t \in \mathcal{M}} \left( \mathbf{u}_{i_t}^\top \mathbf{x}_{i_t,t} - \frac{1}{K} \sum_{j=1}^K \mathbf{u}_j^\top \mathbf{x}_{j,t} \right)^2. \quad (15)$$

In particular, setting  $b = 0$  corresponds to running  $K$  independent (first-order) Perceptron algorithms (no task relatedness), while letting  $b$  go to infinity is optimal only when each

one of the margin deviation terms in (15) is zero (maximal task relatedness). Notice that setting  $b = 1$  gives

$$\begin{aligned} (14) &= \frac{2K-1}{K^2} \sum_{i=1}^K \|\mathbf{u}_i\|^2 \\ &+ \frac{2K-1}{K} \sum_{t \in \mathcal{M}} \left( \mathbf{u}_{i_t}^\top \mathbf{x}_{i_t,t} - \frac{1}{K} \sum_{j=1}^K \mathbf{u}_j^\top \mathbf{x}_{j,t} \right)^2 \end{aligned}$$

yielding a gain of  $K$  whenever the  $K$  tasks are significantly related, as measured by (15).

**Proof of Theorem 6:** Although the general structure of the analysis is based on the second-order Perceptron proof, we use the special properties of  $I + M_t$  in order to compute the contribution of  $K$  and  $b$  to the final bound.

Let  $t = t_s$  be the time step when the  $s$ -th mistake occurs. We write

$$\begin{aligned} \mathbf{v}_s^\top A_s^{-1} \mathbf{v}_s &= (\mathbf{v}_{s-1} + y_t \phi_{i_t,t})^\top A_s^{-1} (\mathbf{v}_{s-1} + y_t \phi_{i_t,t}) \\ &\quad (\text{from the update rule in Figure 3}) \\ &= \mathbf{v}_{s-1}^\top A_s^{-1} \mathbf{v}_{s-1} + 2y_t \mathbf{v}_{s-1}^\top A_s^{-1} \phi_{i_t,t} \\ &\quad + \phi_{i_t,t}^\top A_s^{-1} \phi_{i_t,t} \\ &\leq \mathbf{v}_{s-1}^\top A_s^{-1} \mathbf{v}_{s-1} + \phi_{i_t,t}^\top A_s^{-1} \phi_{i_t,t} \\ &\leq \mathbf{v}_{s-1}^\top A_{s-1}^{-1} \mathbf{v}_{s-1} + \phi_{i_t,t}^\top (I + M_t)^{-1} \phi_{i_t,t}. \end{aligned} \quad (16)$$

In order to prove the first inequality, note that on the  $s$ -th mistaken trial  $A_s = A'_t$  and  $y_t \mathbf{v}_{s-1}^\top (A'_t)^{-1} \phi_{i_t,t} \leq 0$ . In order to prove the second inequality note that  $A_s - A_{s-1}$  and  $A_s - (I + M_t)$  are both positive semidefinite.

We now focus on computing the quadratic form  $\phi_{i_t,t}^\top (I + M_t)^{-1} \phi_{i_t,t}$ . Recall that  $\Phi_t = \sum_{j=1}^K \phi_{j,t}$  is the sum of the orthonormal vectors  $\phi_{1,t}, \dots, \phi_{i_t,t}, \dots, \phi_{K,t}$ . Thus, from the very definition of  $M_t$  in Figure 3 it is easy to verify that

$$(I + M_t) \phi_{i_t,t} = (b(K-1) + 1) \phi_{i_t,t} - \frac{b(K-1)}{K} \Phi_t. \quad (17)$$

Also, since  $M_t \Phi_t = 0$ , we have  $(I + M_t) \Phi_t = \Phi_t$ , and thus  $(I + M_t)^{-1} \Phi_t = \Phi_t$ . Hence (17) allows us to write

$$\phi_{i_t,t} = (b(K-1) + 1) (I + M_t)^{-1} \phi_{i_t,t} - \frac{b(K-1)}{K} \Phi_t.$$

Taking the inner product of both sides with  $\phi_{i_t,t}$ , and solving for  $\phi_{i_t,t}^\top (I + M_t)^{-1} \phi_{i_t,t}$  yields

$$\phi_{i_t,t}^\top (I + M_t)^{-1} \phi_{i_t,t} = \frac{K(b+1) - b}{bK(K-1) + K}.$$

Substituting this into (16), recalling that  $\mathbf{v}_0 = \mathbf{0}$ , and summing over  $s = 1, \dots, m$  we obtain

$$\mathbf{v}_m^\top A_m^{-1} \mathbf{v}_m \leq m \frac{K(b+1) - b}{bK(K-1) + K}.$$

A lower bound on the left-hand side can be obtained in the standard way (details omitted),

$$\sqrt{\mathbf{v}_m^\top A_m^{-1} \mathbf{v}_m} \geq \frac{m - \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u})}{\|A_m^{1/2} \mathbf{u}\|}.$$



Thus we get  $\mathbf{v}_m^\top A_m^{-1} \mathbf{v}_m \geq \frac{(m - \sum_{t \in \mathcal{M}} \ell_t(\mathbf{u}))^2}{\mathbf{u}^\top A_m \mathbf{u}}$ . Finally, recall that by construction

$$\begin{aligned} & \mathbf{u}^\top A_m \mathbf{u} \\ &= \sum_{i=1}^K \|\mathbf{u}_i\|^2 + bK \sum_{t \in \mathcal{M}} \left( \mathbf{u}_{i_t}^\top \mathbf{x}_{i_t, t} - \frac{1}{K} \sum_{j=1}^K \mathbf{u}_j^\top \mathbf{x}_{j, t} \right)^2. \end{aligned}$$

Putting together and solving for  $m$  gives the desired bound.  $\blacksquare$

### 6.1 Implementation in dual variables

Like the second-order multitask Perceptron algorithm, also MMPERC can be formulated in dual variables. Due to the need to handle  $K$  instance vectors at a time, the implementation we sketch below has an extra linear dependence on  $K$ , as compared to the one in Section 4.1.

Let  $t = t_s$  be the trial when the  $s$ -th mistake occurs,  $\mathbf{z}_r$  be the vector  $\mathbf{z}_r = \sqrt{bK}(\phi_{i_r, r} - \Phi_r/K)$ ,  $S_s$  be the matrix whose columns are the vectors  $\phi_{i_r, r}$  corresponding to mistaken time steps  $r$  up to time  $t$ , and  $Z_s$  be the matrix whose columns are the vectors  $\mathbf{z}_r$  corresponding to mistaken time steps  $r$  up to time  $t$ . It is easy to verify that  $A_s = I + Z_s Z_s^\top$  and  $\mathbf{v}_{s-1} = S_s \mathbf{y}_s$  where  $\mathbf{y}_s$  is as in Section 4.1. From the inversion formula (e.g., [17, Ch. 0])

$$(I + Z_s Z_s^\top)^{-1} = I - Z_s (I + Z_s^\top Z_s)^{-1} Z_s^\top$$

we see that the margin  $\mathbf{v}_{s-1}^\top (A_t')^{-1} \phi_{i_t, t}$  in Figure 3 can be computed as

$$\begin{aligned} & \mathbf{v}_{s-1}^\top (A_t')^{-1} \phi_{i_t, t} \\ &= \mathbf{y}_s^\top S_s^\top \phi_{i_t, t} - \mathbf{y}_s^\top S_s^\top Z_s (I + Z_s^\top Z_s)^{-1} Z_s^\top \phi_{i_t, t}. \end{aligned}$$

Calculating the vectors  $S_s^\top \phi_{i_t, t}$  and  $Z_s^\top \phi_{i_t, t}$  in the above expression takes  $\mathcal{O}(s)$  inner products while other  $\mathcal{O}(s)$  inner products are required to incrementally compute  $S_s^\top Z_s$  from  $S_{s-1}^\top Z_{s-1}$ . Finally, when calculating the inverse  $(I + Z_s^\top Z_s)^{-1}$  we exploit the same updating scheme of Section 4.1,

$$I_s + Z_s^\top Z_s = \begin{bmatrix} I_{s-1} + Z_{s-1}^\top Z_{s-1} & Z_{s-1}^\top \mathbf{z}_s \\ \mathbf{z}_s^\top Z_{s-1} & 1 + \mathbf{z}_s^\top \mathbf{z}_s \end{bmatrix},$$

where  $Z_{s-1}^\top \mathbf{z}_s$  and  $\mathbf{z}_s^\top \mathbf{z}_s$  require  $\mathcal{O}(Ks)$  and  $\mathcal{O}(K)$  inner products, respectively. Hence  $(I_s + Z_s^\top Z_s)^{-1}$  can be computed from  $(I_{s-1} + Z_{s-1}^\top Z_{s-1})^{-1}$  with  $\mathcal{O}(Ks)$  extra inner products and  $\mathcal{O}(s^2)$  additional scalar multiplications.

## 7 Spectral co-regularization

An extreme case of multitask learning is the *multiview* setting, where all tasks share the same label. In the multiview protocol, at each time step  $t$  the learner receives  $K$  instances  $\mathbf{x}_{1,t}, \dots, \mathbf{x}_{K,t} \in \mathbb{R}^d$ , predicts with  $\hat{y}_t \in \{-1, 1\}$ , and then receives the correct binary label  $y_t$ , which—unlike the general multitask case—is the same for all instances. What distinguishes multiview learning from a standard online binary classification task, defined on instances of the form  $\mathbf{x}_t = (\mathbf{x}_{1,t}, \dots, \mathbf{x}_{K,t})$ , is that in multiview one postulates the existence of  $K$  vectors  $\mathbf{u}_1, \dots, \mathbf{u}_K$  such that *each*

$\mathbf{u}_i$  is a good linear classifier for the corresponding sequence  $(\mathbf{x}_{i,1}, y_1), (\mathbf{x}_{i,2}, y_2), \dots$  of examples. In this respect a natural baseline for online multiview learning is the algorithm that chooses a random index  $i \in \{1, \dots, K\}$  and then runs a Perceptron algorithm on the sequence of examples  $(\mathbf{x}_{i,t}, y_t)$  for  $t \geq 1$ . Equivalently, we may think of running  $K$  Perceptrons in parallel and then average their mistakes (see the remark after Theorem 10).

In this section we design multiview learning algorithms that in certain cases are able to perform significantly better than the above baseline. In order to do so, we arrange the  $K$   $d$ -dimensional instances  $\mathbf{x}_{1,t}, \dots, \mathbf{x}_{K,t}$  in a  $d \times K$  instance matrix  $X_t$  and penalize diversity among the reference vectors  $\mathbf{u}_1, \dots, \mathbf{u}_K$  using a matrix norm of the  $d \times K$  matrix  $U = [\mathbf{u}_1, \dots, \mathbf{u}_K]$ .

We focus our attention on matrix norms that are unitarily invariant. Such norms are of the form  $\|U\|_f = f(\sigma_U)$ , where  $\sigma_U = (\sigma_1, \dots, \sigma_r)$  is the vector of the ordered singular values  $\sigma_1 \geq \dots \geq \sigma_r \geq 0$  of  $U$  and  $f : \mathbb{R}^r \rightarrow \mathbb{R}$ , with  $r = \min\{K, d\}$ , is an absolutely symmetric function—that is,  $f$  is invariant under coordinate permutations and sign-changes.

Matrix norms of this form control the distribution of the singular values of  $U$ , thus acting as spectral co-regularizers for the reference vectors (see, e.g., [4] for very recent developments on this subject). Known examples are the Schatten  $p$ -norms,  $\|U\|_{s_p} \stackrel{\text{def}}{=} \|\sigma_U\|_p$ . For instance, the Schatten 2-norm is the Frobenius norm. For  $p = 1$  the Schatten  $p$ -norm becomes the trace norm, a good proxy for the rank of  $U$ , since  $\|U\|_{s_1} = \|\sigma_U\|_1 \approx \|\sigma_U\|_0 = \text{rank of } U$ .

In order to obtain a multiview bound that depends on  $\|\sigma_U\|_p$ , we extend the dual norm analysis of Section 5 to matrices. We start by defining the matrix version of the quasi-additive algorithm of [14, 20]. We remark that matrix versions of the EG algorithm and the Winnow algorithm (related to specific instances of the quasi-additive algorithm) have been proposed and analyzed in [28, 29, 30]. When dealing with the trace norm regularizer, their algorithms could be specialized to our multiview framework to obtain mistake bounds comparable to ours. See the brief discussion at the end of this section.

The quasi-additive matrix algorithm maintains a  $d \times K$  matrix  $W$ . Initially,  $W_0$  is the zero matrix. If  $s - 1$  mistakes have been made in the first  $t - 1$  time steps, then the prediction at time  $t$  is  $\text{SGN}(\langle W_{s-1}, X_t \rangle)$ , where  $X_t$  is the  $d \times K$  matrix  $[\mathbf{x}_{1,t}, \dots, \mathbf{x}_{K,t}]$  in which  $\mathbf{x}_{i,t}$  is the instance vector associated with the  $i$ -th view at time  $t$ , and  $\langle W_{s-1}, X_t \rangle$  is the standard matrix inner product  $\langle W_{s-1}, X_t \rangle = \text{TR}(W_{s-1}^\top X_t)$ .

If a mistake occurs at time  $t$ , then  $W_{s-1}$  is updated with  $W_s = \nabla_{\frac{1}{2}} \|V_s\|_f^2$  where, in turn, the  $d \times K$  matrix  $V_s$  is updated using a matrix Perceptron rule,  $V_s = V_{s-1} + y_t X_t$ .

A useful property of norms  $\|U\|_f = f(\sigma_U)$  is that their duals are easily computed.

**Theorem 7 [21, Theorem 2.4]** If  $f$  is absolutely symmetric and  $\|U\|_f = f(\sigma_U)$ , then  $\|U\|_{f^*} = f^*(\sigma_U)$  where  $f^*$  is the convex dual of  $f$ .

In the case of Schatten  $p$ -norms, we have that the dual of vector norm  $\|\cdot\|_p$  is vector norm  $\|\cdot\|_q$ , where  $q = p/(p - 1)$

is the dual coefficient of  $p$ .

An important feature of the quasi-additive algorithms for vectors is that the mapping  $\mu : \mathbf{v} \mapsto \mu(\mathbf{v}) = \nabla \frac{1}{2} \|\mathbf{v}\|^2$  is invertible whenever the vector norm  $\|\cdot\|$  satisfies certain regularity properties (see, e.g., [8, page 294]). We call such norms *Legendre*. Hence, we “do not lose information” when the primal weight vector  $\mathbf{v}$  is mapped to the weight vector  $\mathbf{w} = \mu(\mathbf{v})$  used for prediction. In particular, we always have that  $\mu^{-1}(\mathbf{v}) = \nabla \frac{1}{2} \|\mathbf{v}\|_*^2$ , where  $\|\cdot\|_*$  is the dual norm of a Legendre norm  $\|\cdot\|$  (see, e.g., [8, Lemma 11.5]). This property is preserved when the algorithm is applied to matrices. This is shown by the following result where, without loss of generality, we prove the property for  $f(\sigma_U)$  rather than  $\frac{1}{2}f(\sigma_U)^2$ . (In fact, if  $f$  is Legendre, then  $\frac{1}{2}f^2$  is also Legendre).

**Theorem 8** *Let  $f$  be a Legendre function. If  $\|U\|_f = f(\sigma_U)$  then  $(\nabla \|\cdot\|_f)^{-1} = \nabla \|\cdot\|_{f^*}$ .*

The following result will be useful.

**Theorem 9 [21, Theorem 3.1]** *Let  $U \text{DIAG}[\sigma_A] V^\top$  be an SVD decomposition of a matrix  $A$ . If  $\|\cdot\|_f$  is a matrix norm such that  $\|A\|_f = f(\sigma_A)$  for  $f$  Legendre, then  $\nabla f(\sigma_A) = \sigma_{\nabla \|A\|_f}$ . Moreover,  $\nabla \|A\|_f = U \text{DIAG}[\nabla f(\sigma_A)] V^\top$ .*

**Proof of Theorem 8:** If  $A = U \text{DIAG}[\sigma_A] V^\top$ , then by Theorem 9

$$\nabla \|A\|_f = U \text{DIAG}[\nabla f(\sigma_A)] V^\top = U \text{DIAG}[\sigma_{\nabla \|A\|_f}] V^\top.$$

Therefore, using Theorem 9,

$$\begin{aligned} \nabla \left\| (\nabla \|A\|_f) \right\|_{f^*} &= U \text{DIAG} \left[ \nabla f^*(\nabla f(\sigma_A)) \right] V^\top \\ &= U \text{DIAG}[\sigma_A] V^\top \quad (f \text{ is Legendre}) \\ &= A \end{aligned}$$

concluding the proof.  $\blacksquare$

We now develop a general analysis of the quasi-additive matrix algorithms, and then specialize it (in Theorem 10 below) to a multiview algorithm operating with a Schatten  $p$ -norm regularizer.

We start by adapting the dual norm proof of Section 5 to an arbitrary matrix norm  $\|A\|_f = f(\sigma_A)$ , where  $f$  is Legendre. Let  $V_m$  be the primal weight matrix after any number  $m$  of mistakes. By Taylor expanding  $\frac{1}{2} \|V_s\|_f^2$  around  $V_{s-1}$  for each  $s = 1, \dots, m$ , and using  $y_t \langle W_{s-1}, X_t \rangle \leq 0$ , we get

$$\frac{1}{2} \|V_m\|_f^2 \leq \sum_{s=1}^m D(V_s \| V_{s-1})$$

where  $D(V_s \| V_{s-1})$  is the matrix Bregman divergence  $\frac{1}{2} (\|V_s\|_f^2 - \|V_{s-1}\|_f^2) - y_t \langle W_{s-1}, X_t \rangle$ .

Fix any  $d \times K$  matrix  $U$ . First, we derive a matrix version of the convex inequality for vector norms. We use a classical result by von Neumann (see, e.g., [18, p. 182]) stating that  $\langle V, U \rangle \leq \sigma_V^\top \sigma_U$  for any two  $d \times K$  matrices  $U$  and  $V$ . We have

$$\begin{aligned} \|V_m\|_f \|U\|_{f^*} &= f(\sigma_{V_m}) f^*(\sigma_U) \quad (\text{by Theorem 7}) \\ &\geq \sigma_{V_m}^\top \sigma_U \quad (\text{by the convex ineq. for norms}) \\ &\geq \langle V_m, U \rangle \quad (\text{by von Neumann's ineq.}). \end{aligned}$$

In addition, we have  $\langle U, V_s \rangle = \langle U, V_{s-1} \rangle + y_t \langle U, X_t \rangle$ . Thus we obtain

$$\|V_m\| \geq \frac{\langle U, V_m \rangle}{\|U\|_{f^*}} \geq \frac{Km - \sum_t \ell_t(U)}{\|U\|_{f^*}}$$

where  $\ell_t(U) \stackrel{\text{def}}{=} \sum_{i=1}^K [1 - y_t \mathbf{u}_i^\top \mathbf{x}_{i,t}]_+ = \sum_{i=1}^K \ell_t(\mathbf{u}_i)$ . Solving for  $m$  gives

$$m \leq \frac{1}{K} \sum_{t \in \mathcal{M}} \ell_t(U) + \frac{\|U\|_{f^*}}{K} \sqrt{2 \sum_{s=1}^m D(V_s \| V_{s-1})}. \quad (18)$$

Equation (18) is our general starting point for analyzing multiview algorithms working under spectral co-regularization. The analysis reduces to bounding from above the second-order term  $D(\cdot \| \cdot)$  of the specific matrix norm  $\|\cdot\|_f$  under consideration.

For the rest of this section we focus on the Schatten  $2p$ -norm  $\|V\|_{s_{2p}} = \|\sigma_V\|_{2p}$ , where  $V$  is a generic  $d \times K$  matrix, and  $p$  is a positive integer (thus  $2p$  is an even number  $\geq 2$ ). Note that, in general,  $\|V\|_{s_{2p}}^2 = \text{TR}((V^\top V)^p)^{1/p}$ .

In order to prove our main result, stated below, we use some facts from differential matrix calculus. A standard reference on this subject is [23], to which the reader is referred.

**Theorem 10** *The number of mistakes  $m$  made by the  $2p$ -norm matrix Perceptron, run on any sequence  $(X_1, y_1), (X_2, y_2), \dots$  satisfies, for any  $d \times K$  matrix  $U$ ,*

$$\begin{aligned} m \leq \frac{1}{K} \sum_{t \in \mathcal{M}} \ell_t(U) &+ (2p-1) \left( \frac{X_{s_{2p}} \|U\|_{s_{2q}}}{K} \right)^2 \\ &+ \frac{X_{s_{2p}} \|U\|_{s_{2q}}}{K} \sqrt{\frac{2p-1}{K} \sum_{t \in \mathcal{M}} \ell_t(U)} \end{aligned}$$

where  $X_{s_{2p}} = \max_{t \in \mathcal{M}} \|X_t\|_{s_{2p}}$ ,  $\|U\|_{s_{2q}}$  is the Schatten  $2q$ -norm of  $U$ , with  $2q = \frac{2p}{2p-1}$ , and  $\mathcal{M}$  is the set of mistaken trial indices.

**Remark** Similarly to the vector case, when the parameter  $p$  is chosen to be logarithmic in  $r = \min\{d, K\}$ , the  $p$ -norm matrix Perceptron penalizes diversity using the trace norm of  $U$ . If the vectors  $\mathbf{u}_i$  span a subspace of size  $\ll K$ , and instances tend to have  $K$  nonzero singular values of roughly the same magnitude, then  $\|U\|_{s_{2q}} \approx \|U\|_{s_2}$  while  $X_{s_{2p}}^2 \approx X_{s_\infty}^2 \approx X_{s_2}^2/K$ . Hence this choice of  $p$  leads (at least in the linearly separable case) to a factor  $K$  improvement over the bound achieved by the matrix algorithm based on the Frobenius norm ( $p = 1$  in Theorem 10), which amounts to running  $K$  independent Perceptrons in parallel and then average their mistakes.

The following trace inequality is our main technical lemma.

**Lemma 11** *Let  $A, B$  be positive semidefinite matrices, of size  $d \times d$  and  $K \times K$  respectively, with the same nonzero eigenvalues. Let  $X$  be an arbitrary real matrix of size  $d \times K$ . Then, for any pair on nonnegative exponents  $l, g \geq 0$ , we have  $\text{TR}(X^\top A^l X B^g) \leq (\text{TR}(X^\top X)^p)^{1/p} (\text{TR} A^{(l+g)q})^{1/q}$  where  $1/p + 1/q = 1$ ,  $p \geq 1$ .*

**Proof of Lemma 11:** We first consider the case  $l \leq g$ . By the Cauchy-Schwartz and Holder's inequalities applied to traces [23, Ch.11] we have

$$\begin{aligned} \text{TR}(X^\top A^l X B^g) &= \text{TR}(B^{(g-l)/2} X^\top A^l X B^{(g+l)/2}) \quad (19) \\ &\leq \text{TR}(X^\top A^{2l} X B^{g-l})^{1/2} \text{TR}(X^\top X B^{g+l})^{1/2} \\ &\leq \text{TR}(X^\top A^{2l} X B^{g-l})^{1/2} T_p(X^\top X)^{1/2} T_q(B^{g+l})^{1/2} \end{aligned}$$

where we used the shorthand  $T_r(Z) = (\text{TR} Z^r)^{1/r}$ . In the case when  $l > g$  we can simply swap the matrices  $X^\top A^l$  and  $X B^g$  and reduce to the previous case.

We now recursively apply the above argument to the left-hand side of (19). Recalling that  $T_q(A) = T_q(B)$  and  $T_p(X^\top X) = T_p(X X^\top)$ , after  $n$  steps we obtain

$$\begin{aligned} \text{TR}(X^\top A^l X B^g) &\leq (\text{TR}(X^\top A^{l'} X B^{g'}))^{1/2^n} \times \\ &\times T_p(X^\top X)^{\sum_{i=1}^n (1/2)^i} T_q(B^{g+l})^{\sum_{i=1}^n (1/2)^i} \end{aligned}$$

for some pair of exponents  $l', g' \geq 0$  such that  $l' + g' = l + g$ . Since for any such pair  $l', g'$ , we have  $\text{TR}(X^\top A^{l'} X B^{g'}) < \infty$ , we can take the limit as  $n \rightarrow \infty$ . Recalling that  $\sum_{i=1}^\infty (1/2)^i = 1$  completes the proof. ■

**Proof of Theorem 10:** We set for brevity  $G : \mathbb{R}^{d \times K} \rightarrow \mathbb{R}$ ,

$$G(V) = \frac{1}{2} \text{TR}((V^\top V)^p)^{1/p} = \frac{1}{2} \|V\|_{s_{2p}}^2.$$

Thus in our case

$$D(V_s \| V_{s-1}) = G(V_{s-1} + y_t X_t) - G(V_{s-1}) - y_t \langle W_{s-1}, X_t \rangle.$$

Since  $G(V)$  is twice<sup>4</sup> continuously differentiable, by the mean-value theorem we can write

$$D(V_s \| V_{s-1}) = \frac{1}{2} \text{VEC}(X_t)^\top H_G(\xi) \text{VEC}(X_t), \quad (20)$$

where  $\text{VEC}(X)$  is the standard columnwise vectorization of a matrix  $X$ ,  $H_G$  denotes the Hessian matrix of (matrix) function  $G$  and  $\xi$  is some matrix on the line connecting  $V_{s-1}$  to  $V_s$ . Using the rules of matrix differentiation, the gradient  $\nabla G$  of  $G$  is  $\nabla G(V) = c(V) \text{VEC}(D)^\top$  where we set for brevity  $D = V B^{p-1}$ ,  $c(V) = \text{TR}(B^p)^{1/p-1}$ , with  $B = V^\top V$ . Taking the second derivative  $H_G = \nabla \nabla G$  gives  $H_G(V) = \text{VEC}(D) \nabla(c(V)) + c(V) \nabla(D)$ . Now, recalling the definition of  $c(V)$ , it is not hard to show that  $\text{VEC}(D) \nabla(c(V))$  is the  $Kd \times Kd$  matrix

$$2(1-p) \text{TR}(B^p)^{1/p-2} \text{VEC}(D) \text{VEC}(D)^\top.$$

Since  $p \geq 1$  this matrix is negative semidefinite, and we can disregard it when bounding from the above the quadratic form (20). Thus we continue by considering only the second term  $c(V) \nabla(D)$  of the Hessian matrix. We have

$$\nabla(D) = (B^{p-1} \otimes I_d) + (I_k \otimes V) \nabla(B^{p-1}),$$

<sup>4</sup>In fact  $G$  is  $C^\infty$  everywhere but (possibly) in zero, since  $\text{TR}((V^\top V)^p)$  is just a polynomial function of the entries of  $V$ . Moreover  $\text{TR}((V^\top V)^p) = 0$  if and only if  $V$  is the zero matrix.

where

$$\nabla(B^{p-1}) = \left( \sum_{\ell=0}^{p-2} B^\ell \otimes B^{p-2-\ell} \right) (I_{K^2} + T_K) (I_k \otimes V^\top),$$

and  $T_K$  is the  $K^2 \times K^2$  commutation matrix such that  $T_K \text{VEC}(M) = \text{VEC}(M^\top)$  for any  $K \times K$  matrix  $M$ . Putting together

$$\begin{aligned} (20) &\leq \frac{c(V)}{2} \text{VEC}(X_t)^\top (B^{p-1} \otimes I_d) \text{VEC}(X_t) \\ &\quad + \frac{c(V)}{2} \text{VEC}(X_t)^\top (I_k \otimes V) \Sigma \times \\ &\quad \times (I_{K^2} + T_K) (I_k \otimes V^\top) \text{VEC}(X_t), \quad (21) \end{aligned}$$

where we used the shorthand  $\Sigma = \sum_{\ell=0}^{p-2} B^\ell \otimes B^{p-2-\ell}$ . We now bound the two terms in the right-hand side of (21). By well-known relationships between Kronecker products and the  $\text{VEC}$  operator (see [23, Ch. 3]) we can write

$$\begin{aligned} \frac{c(V)}{2} \text{VEC}(X_t)^\top (B^{p-1} \otimes I_d) \text{VEC}(X_t) \\ = \frac{c(V)}{2} \text{TR}(X_t^\top X_t B^{p-1}) \leq \frac{1}{2} (\text{TR}(X_t^\top X_t))^p, \end{aligned}$$

independent of  $V$ . The majorization follows from Holder's inequality applied to the positive semidefinite matrices  $X_t^\top X_t$  and  $B^{p-1}$ . Moreover, it is easy to see that the symmetric matrices  $\Sigma$  and  $T_K$  commute, thereby sharing the same eigenspace. Hence,  $\Sigma(I_{K^2} + T_K) \preceq 2\Sigma$ , and we can bound from above the second term in (21) by

$$c(V) \text{VEC}(X_t)^\top \sum_{\ell=0}^{p-2} B^\ell \otimes A^{p-1-\ell} \text{VEC}(X_t),$$

where we set  $A = V V^\top$ . Again, [23, Ch. 3] allows us to rewrite this quadratic form as the sum of traces

$$c(V) \sum_{\ell=0}^{p-2} \text{TR}(X_t^\top A^{p-1-\ell} X_t B^\ell).$$

Since  $A$  and  $B$  have the same nonzero eigenvalues, we can apply Lemma 11 to each term and put together as in (21). After simplifying we get

$$(20) \leq \frac{1}{2} (2p-1) (\text{TR}(X_t^\top X_t))^p = \frac{1}{2} (2p-1) \|X_t\|_{s_{2p}}^2.$$

The desired bound is then obtained by plugging back into (18), solving the resulting inequality for  $m$ , and overapproximating. ■

The result of Theorem 10 is similar to those obtained in [28, 29, 30]. However, unlike these previous results, our matrix algorithm has no learning rate to tune (a property inherited from the vector  $p$ -norm Perceptron of [13]) and works for arbitrary nonsquare matrices  $U$ . We also observe that the prediction  $\hat{y}_t = \text{SGN}(\text{TR}(W_{s-1}^\top X_t))$  of the  $p$ -norm matrix Perceptron reduces to computing the sign of  $\text{TR}((V_{s-1}^\top V_{s-1})^{p-1} V_{s-1}^\top X_t)$  (recall the expression for  $\nabla G$  calculated in the proof of Theorem 10). Since matrix  $V_s$  is updated additively, it is clear that both  $V_{s-1}^\top V_{s-1}$  and  $V_{s-1}^\top X_t$  do depend on instance vectors  $x_{i,t}$  only through inner products. This allows us to turn our  $p$ -norm matrix Perceptron into a kernel-based algorithm, and repeat the analysis given here using a standard RKHS formalism.

## 8 Conclusions and ongoing research

In this work we have studied the problem of learning multiple tasks online using various approaches to formalize the notion of task relatedness.

Our results can be extended in different directions. First, in Sections 3.2 and 6 it might be interesting to devise methods for dynamically adapting the  $b$  parameter as new data are revealed. Second, the mistakes of the second-order algorithm MMPERC have been bounded using a first-order analysis. A more refined analysis should reveal in the bound an explicit dependence on the spectral properties of the data. It is also worth noting that the significance of the mistake bound obtained for MMPERC relies on the fact that the algorithm assumes the tasks to be different, although somewhat related. In the case when the  $K$  observed instances share the same label at each time step (like in multiview learning), we could not devise an algorithm with a significant advantage over the following trivial baseline: run  $K$  Perceptrons in parallel and use the sum of margins to predict. Third, it would be interesting to study the problem of learning multiple tasks when  $K$  predictions have to be output in each step. In this case the main difficulty appears to be the control of the interaction among instances at each time step. Fourth, it would be also interesting to prove *lower* bounds on the number of mistakes, as a function of task relatedness. Finally, since multitask learning problems arise naturally in a variety of settings, spanning from biology to news processing, we plan to complement the theoretical analysis presented in this paper with experimental results, so as to evaluate the empirical performance of our algorithms in real-case scenarios.

**Acknowledgments.** Thanks to Sham Kakade, Massi Pontil, and Francis Bach for useful discussions. We also thank the COLT 2008 reviewers for their comments. This work was supported in part by the PASCAL2 Network of Excellence under EC grant no. 216886. This publication only reflects the authors' views.

## References

- [1] J. ABERNETHY, P.L. BARTLETT & A. RAKHLIN, Multitask learning with expert advice, Proc. 20th COLT, pp. 484–498, Springer, 2007.
- [2] R.K. ANDO & T. ZHANG, A framework for learning predictive structures from multiple tasks and unlabeled data, *JMLR*, 6, pp. 1817–1853, MIT Press, 2005.
- [3] A. ARGYRIOU, T. EVGENIOU & M. PONTIL Multi-Task feature learning, NIPS 19, pp. 41–48, MIT Press, 2007.
- [4] A. ARGYRIOU, C.A. MICCHELLI, M. PONTIL & Y. YING, A spectral regularization framework for multi-task structure learning, NIPS 20, MIT Press, 2008.
- [5] K. AZOURY AND M. WARMUTH, Relative loss bounds for on-line density estimation with the exponential family of distributions, *Machine Learning*, 43, pp. 211–246, 2001.
- [6] U. BREFELD, T. GAERTNER, T. SCHEFFER, & S. WROBEL, Efficient co-regularised least squares regression, Proc. 23rd ICML, 2006.
- [7] N. CESA-BIANCHI, A. CONCONI & C. GENTILE, A second-order Perceptron algorithm, *SIAM Journal on Computing*, 34/3, pp. 640–668, 2005.
- [8] N. CESA-BIANCHI & G. LUGOSI, *Prediction, Learning, and Games*. Cambridge University Press, 2006.
- [9] O. DEKEL, P.M. LONG & Y. SINGER, Online learning of multiple tasks with a shared loss, *JMLR*, 8, pp. 2233–2264, 2007.
- [10] T. EVGENIOU, M. PONTIL & T. POGGIO, Regularization networks and Support Vector Machines, *Advances in Computational Mathematics*, 13/1, pp. 1–50, Springer, 2000.
- [11] T. EVGENIOU, C. MICCHELLI & M. PONTIL, Learning Multiple tasks with kernel methods, *JMLR*, 6, pp. 615–637, MIT Press, 2005.
- [12] Y. FREUND & R.E. SCHAPIRE, Large margin classification using the Perceptron algorithm. *Machine Learning*, 37:3, pp. 277–296, 1999.
- [13] C. GENTILE, The robustness of the  $p$ -norm algorithms, *Machine Learning*, 53, pp. 265–299, 2003.
- [14] A. GROVE, N. LITTLESTONE & D. SCHUURMANS, General convergence results for linear discriminant updates, *Machine Learning*, 43, pp. 173–210, 2001.
- [15] M. HERBSTER, M. PONTIL & L. WAINER, Online learning over graphs, Proc. 22nd ICML, pp. 305–312, ACM Press, 2005.
- [16] L. HOGGEN, *Handbook of Linear Algebra*, Discrete Mathematics and Its Applications, 39, CRC Press, 2006.
- [17] R.A. HORN & C.R. JOHNSON, *Matrix Analysis*. Cambridge University Press, 1985.
- [18] R.A. HORN & C.R. JOHNSON, *Topics in Matrix Analysis*. Cambridge University Press, 1991.
- [19] A. JAGOTA & M.K. WARMUTH, Continuous and discrete time nonlinear gradient descent: relative loss bounds and convergence, Electr. Proc. 5th International Symposium on Artificial Intelligence and Mathematics, 1998. Electronic, <http://rutcor.rutgers.edu/~amai>.
- [20] J. KIVINEN & M. WARMUTH, Relative loss bounds for multidimensional regression problems, *Machine Learning*, 45, pp. 301–329, 2001.
- [21] A.S. LEWIS, The convex analysis of unitarily invariant matrix functions, *Journal of Convex Analysis*, 2, pp. 173–183, 1995.
- [22] N. LITTLESTONE, *Mistake bounds and logarithmic linear-threshold learning algorithms*. Ph.D. Thesis, University of California at Santa Cruz, 1989.
- [23] J.R. MAGNUS, & H. NEUDECKER, *Matrix Differential Calculus with Applications in Statistics and Econometrics*, revised edition. John Wiley, 1999.
- [24] A. MAURER, Bounds for linear multi-task learning, *JMLR*, 7, pp. 117–139, MIT Press, 2006.
- [25] R.T. ROCKAFELLAR, *Convex Analysis*. Princeton University Press, 1970.
- [26] D. ROSENBERG & P. BARTLETT, Rademacher complexity of co-regularized kernel classes, Proc. Artificial Intelligence and Statistics, 2007.
- [27] V. SINDHWANI, P. NIYOGI & M. BELKIN, A co-regularized approach to semi-supervised learning. Proc. ICML Workshop on Learning with Multiple Views, 2005.
- [28] K. TSUDA, G. RAETSCH & M.K. WARMUTH, Matrix exponentiated gradient updates for on-line learning and Bregman projection, *JMLR*, 6, pp. 995–1018, 2005.
- [29] M.K. WARMUTH & D. KUZMIN, Online variance minimization, Proc. 19th COLT, Springer, 2006.
- [30] M.K. WARMUTH, Winoing subspaces. Proc. 24th ICML, pp. 999–1006, ACM Press, 2007.